

# Gene-Environment Interaction Analysis Using Graphic Cards

Daniel Berglund

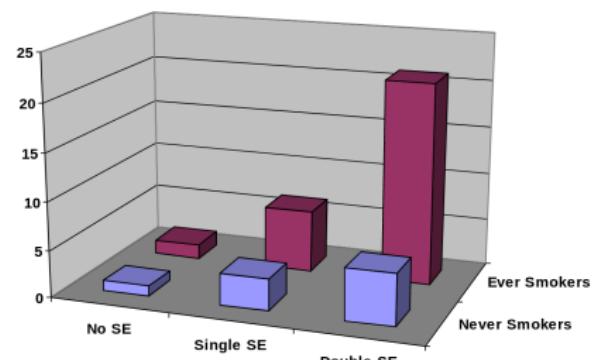
17 February 2015

*Supervisor:*  
Lilit AXNER

*Project provider:*  
Henrik KÄLLBERG

# Problem

- Genetic and environmental factors are known to affect the risks of diseases
- Interaction can exist between these factors



---

Klareskog L, Stolt P, Lundberg K, et. al. A new model for an etiology of rheumatoid arthritis: smoking may trigger HLA-DR (shared epitope)-restricted immune reactions to autoantigens modified by citrullination. *Arthritis Rheum* 2006

# Aim

- Data amounts are increasing, need for better programs
- GPUs have previously shown good results for gene-gene interaction
- Build a program for gene-environment interaction using GPUs based on older programs

# The Data

- 100 000+ genetic markers
- One environment factor, for instance smoking
- Covariates, variables suspected or known to effect the risk, for instance age

# Contingency Table

Genetic marker	Case			Control		
	No SE	1 SE	2 SE	No SE	1 SE	2 SE
Ever smoker	58	192	126	184	146	31
Never smoker	20	72	36	87	104	31

# Relative Risk, Odds and Odds Ratio

$$\text{Relative risk} = \frac{P(Y = 1|X = x_1)}{P(Y = 1|X = x_2)}$$

$$\text{Odds} = \frac{P(Y = 1|X = x)}{1 - P(Y = 1|X = x)}$$

$$\text{Odds ratio} = \frac{\text{Odds}_1}{\text{Odds}_2}$$

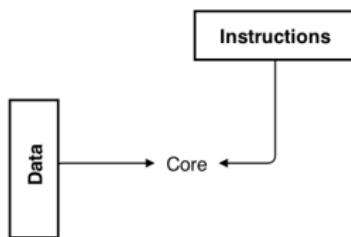
# Logistic Regression

$$\log(odds) = \alpha + \beta X$$

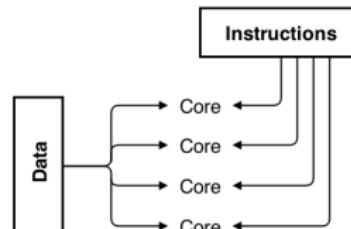
$$Odds\ ratio = e^{\beta}$$

$\beta$  coefficients can be found by using Newtons method and maximum likelihood

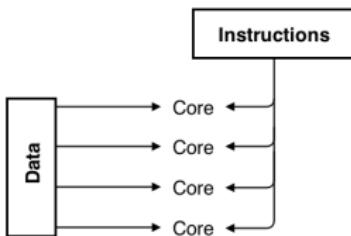
# Concurrency Architectures



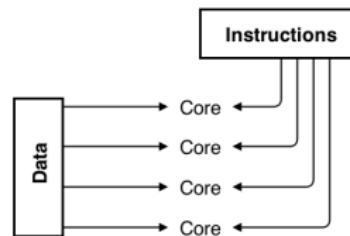
SISD



MISD



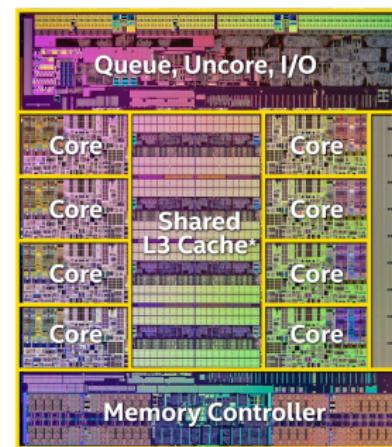
SIMD



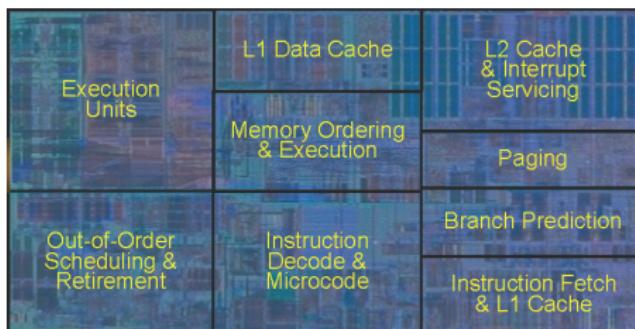
MIMD

# CPU Architecture

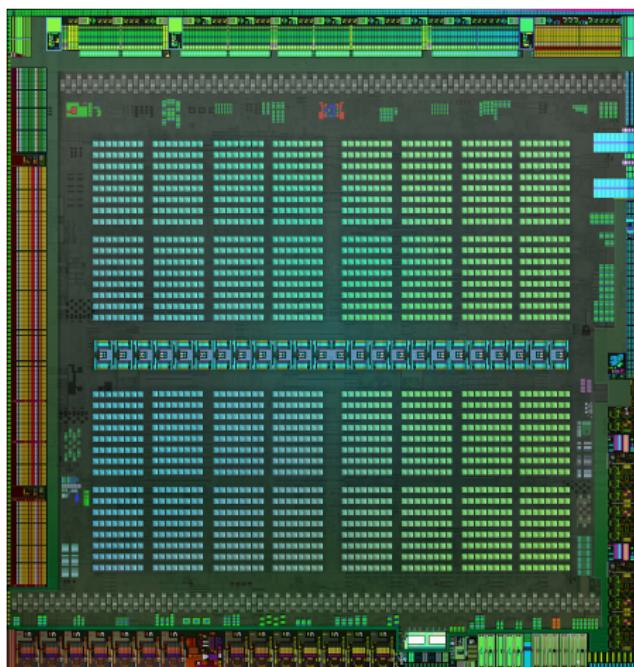
- MIMD
- 8 cores
- Similar performance with single and double precision
- Lots of various optimizations



# CPU Architecture



# GPU Architecture



# GPU Architecture

- Slow but many cores, up to 5000
- Fewer optimizations than CPU
- Separate memory from the normal RAM
- Data has to be explicitly transferred to and from the GPU
- Much better performance with single precision than double
- Single instruction, multiple thread(SIMT)

# CPU versus GPU

- Pick CPU or GPU based on the problem
- CPU for complex control flows
- GPUs need embarrassingly parallelism
- CPUs are becoming more like GPUs and vice versa
- Xeon Phi, CPU accelerators

# CUDA

- Compute Unified Device Architecture
- For NVIDIA GPUs by NVIDIA
- Provides a compiler for kernels
- Supported by various libraries, for instance CUBLAS, Thrust
- Calculations done by kernels

# Kernel Example

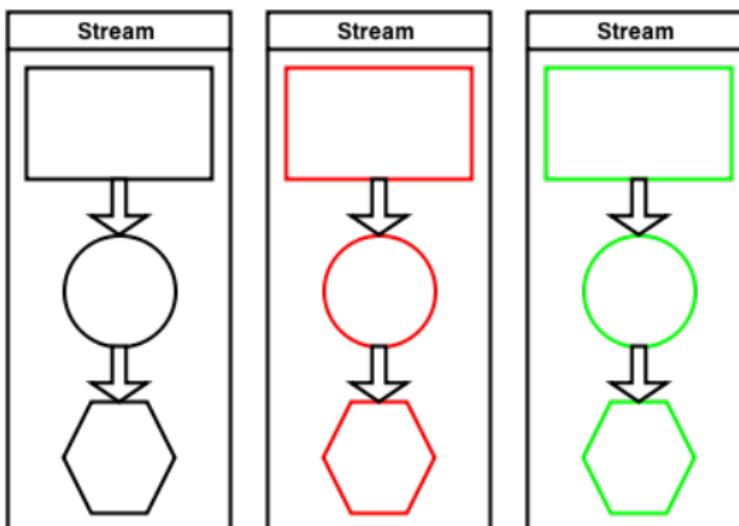
```
__global__ void Add(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}
```

Add<<< N,M >>>(A, B, C);

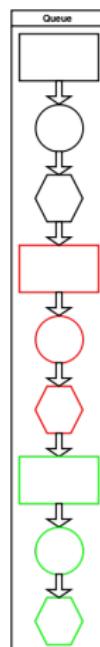
# Streams

- Kernels and transfers can be performed on streams
- Transfers and kernels from different streams can overlap
- Transfers and kernels uses different parts of the GPU

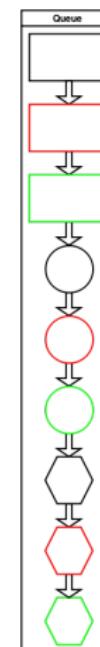
# Streams



# Streams



*Looped over stream*



*Looped over kernel*

# Agile Development

- Code in small units, commonly classes
- Unit testing, test the units in isolation
- Mocks, i.e. fake objects
- If a unit is working incorrectly only its corresponding unit test should fail
- Integration tests can be used to test if the units work properly together

# Wrappers

- Wrappers are used to “fix” interfaces
- Does not perform the task, delegates it

# Wrappers

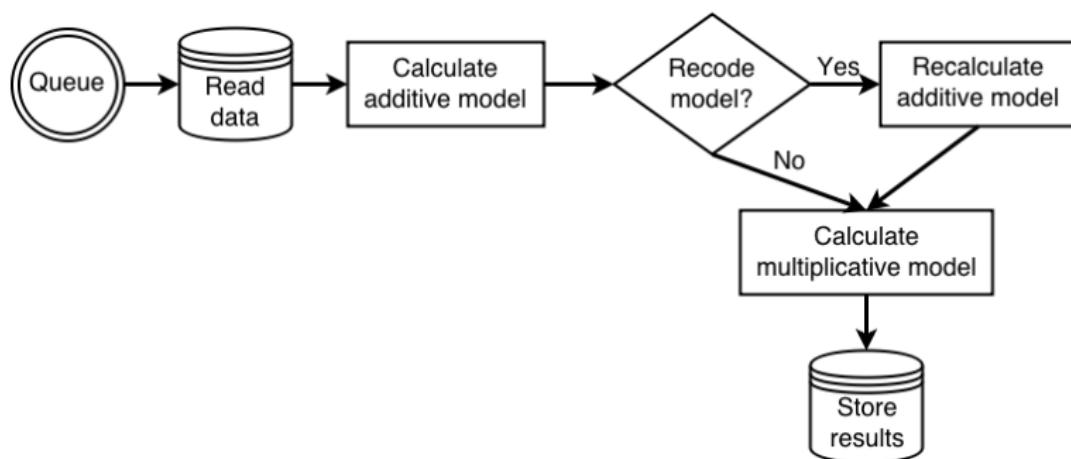
```
cublasSgemv(cublasHandle, transpose, m, n, α, matrix,  
ld_matrix, vector, inc_vector, β, result, inc_result);
```

```
matrixVectorMultiply(matrix, vector, result, α, β);
```

# Structure of Program

- Called CuEira
- Written in C++
- Libraries used are Google Test/Mock, Boost, CUDA, CUBLAS
- Modularised
- Almost full unit test coverage
- Some parts are done using CPU
- Handles each gene environment combination independently
- A number of CPU threads perform the work and each use a GPU stream for calculations

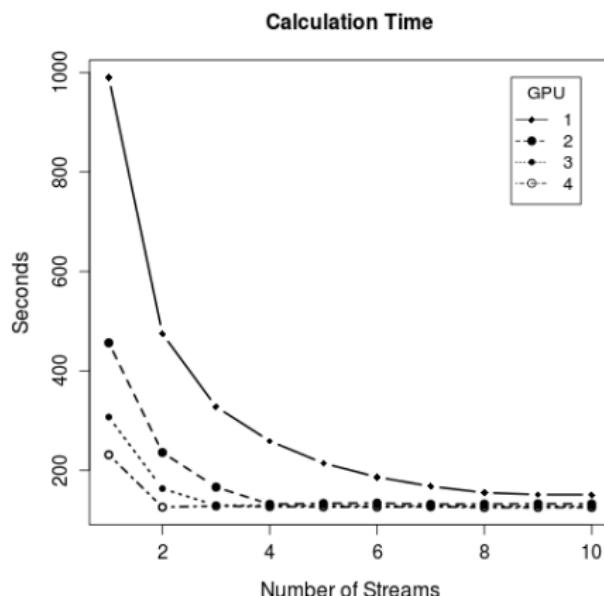
# Overview



# Data and Expectations

- Simulated data
- 10 000 genetic markers
- 2000, 10 000, 100 000, 200 000 individuals
- 0, 5, 10, 20 covariates
- Expected linear efficiency, i.e. good scaling to multiple GPUs

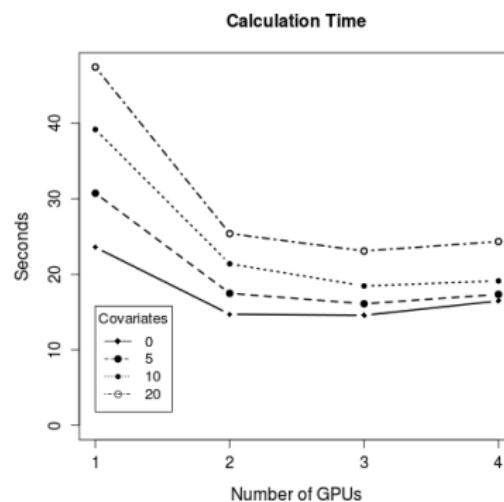
# Streams



*100 000 individuals, 0 covariates*

# Saturated Streams

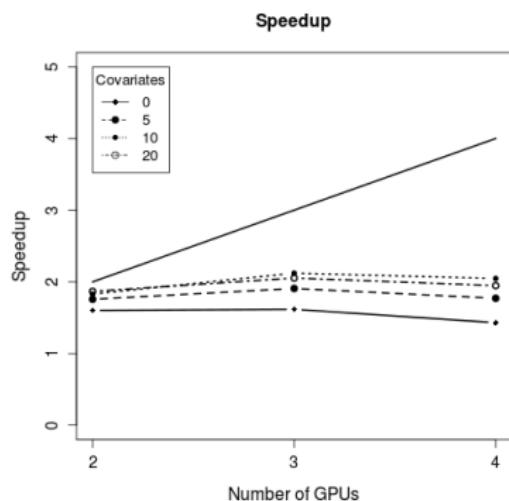
GPUs	Streams
1	9
2	4
3	3
4	2



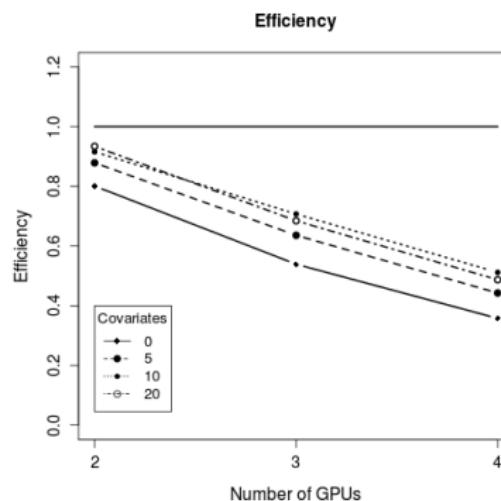
*2000 individuals*

# Speedup and Efficiency

## Versus one GPU

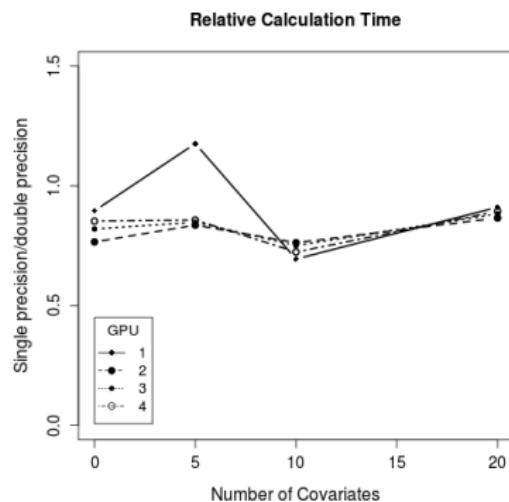


*2000 individuals*

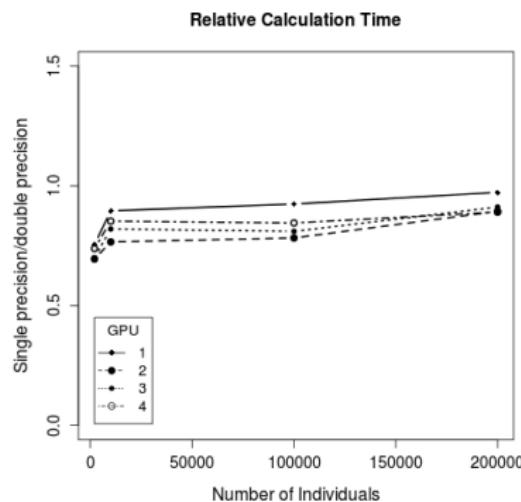


*2000 individuals*

# Single Versus Double Precision

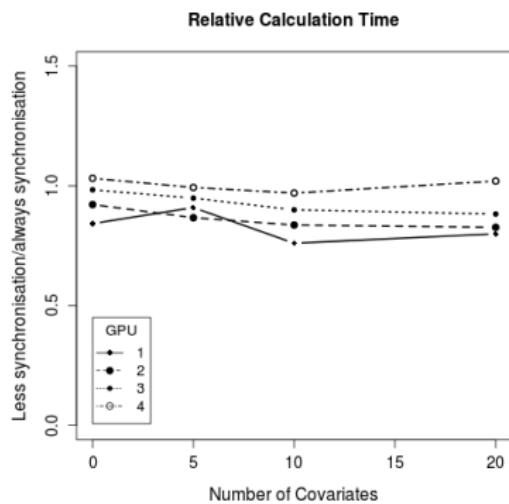


*10000 individuals*

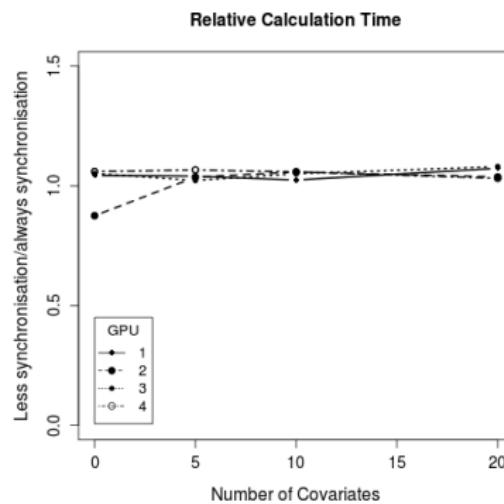


*0 covariates*

# Synchronisation

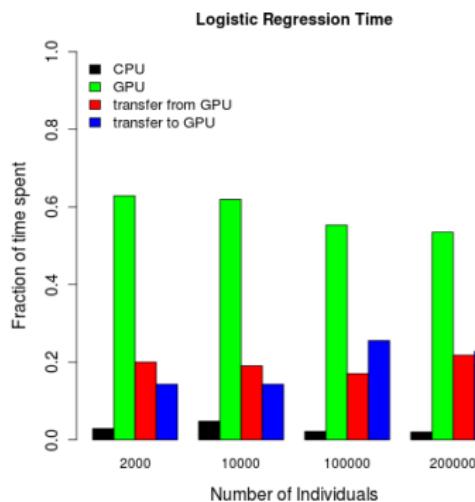


*2000 individuals*

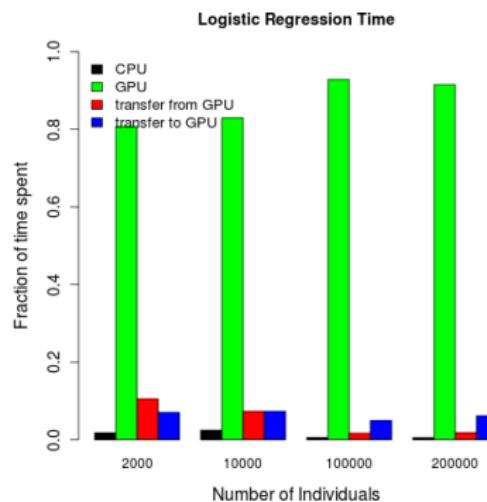


*200 000 individuals*

# Synchronisation

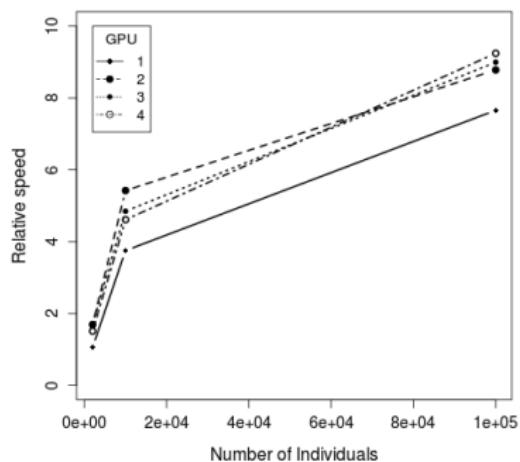


*Less synchronisation*

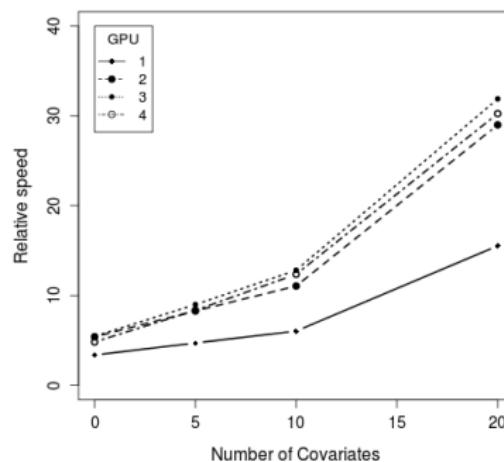


*Always synchronise*

# Comparison with GEISA



0 covariates



10 000 individuals

# Conclusions

- GPUs suits well for interaction, depending on method
- Use one, in some cases two, GPUs with the program
- Single precision gives better performance than double precision
- Increased synchronisation increases performance slightly when number of individuals is above 100 000

# Outlook

- Implement a method for validation of the models
- Move more parts to the GPU to improve the scaling
- Use clusters for more speed
- There is a need for better statistics and definitions of interaction for non binary factors
- For binary factors look into gene-gene interaction methods if further speed is needed

Thanks to

Lilit Axner, Henrik Källberg, KIRC, Michael Schliephake, PDC,  
IMM, Jens Lagergren

CuEira is available at [github.com/Berjiz/CuEira](https://github.com/Berjiz/CuEira)

# Time Distribution for Kernels

\* is element by element multiplication

Covariates Individuals		0 10k	0 100k	20 10k	20 100k
$M1^T \cdot V1$	CUBLAS	66.6	72.2	43	53.9
$M1 \cdot M2$	CUBLAS	18.9	14.4	36.1	25.8
$M1 \cdot V1$	CUBLAS	7.1	8.1	5.9	7.4
$V1 * V2$	Custom	2.9	2.3	11.8	10.6
$\frac{e^{V1}}{1+e^{V1}}$	Custom	1.1	0.9	0.7	0.7
$V1 * \log(V2) +$	Custom	0.8	0.6	0.6	0.5
$(1 - V1) * \log(1 - V2)$	Custom	0.7	0.6	0.5	0.4
$V1 - V2$	Custom	0.7	0.5	0.4	0.4
$V1 * (1 - V1)$	Custom	0.7	0.5	0.4	0.4
<i>Dot product</i>	CUBLAS	0.7	0.3	0.3	0.2
$\Sigma V1$	CUBLAS	0.5	0.1	0.2	0.1

# CPU versus GPU

	CPU	MIC	GPU
Example	Intel i7-5960X	Xeon Phi 7120A	K80
Cores	8	61	4992
Memory	0.1-1 TB	16GB	24 GB
Caches	20 MB	30.5 MB	-
Concurrency	MIMD	MIMD	SIMT
FLOPS	0.4 T	3 T	3, 8.7 T
Price	\$1000	\$4200	\$5000

# Matrix Decomposition

- Pseudo inverse,  $A^+$ , is defined for general matrices
- Can be found by using singular value decomposition

$$A = U\Sigma V^T$$

$$A^+ = V\Sigma^+ U^T$$

# Recoding

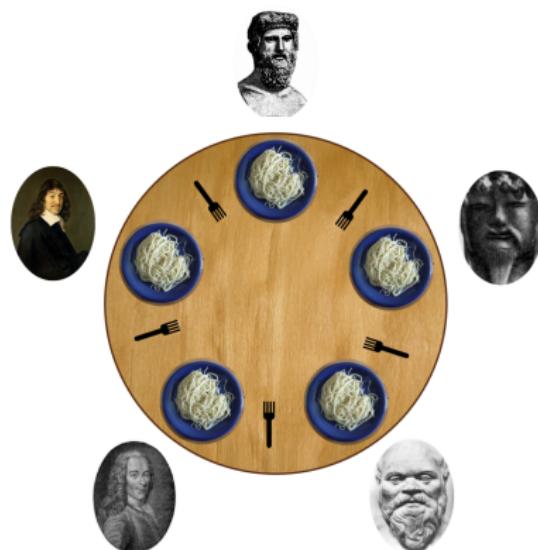
- The measures for additive interaction are defined for positive odds ratios
- Can be adjusted by recoding
- Recoding switches the reference group with the group with lowest risks
- Guarantees that  $OR \geq 1$

# Speedup and Efficiency

$$S(p) = \frac{T(1)}{T(p)}$$

$$E(p) = \frac{S(p)}{p} = \frac{T(1)}{pT(p)}$$

# Concurrency and Dinning Philosophers



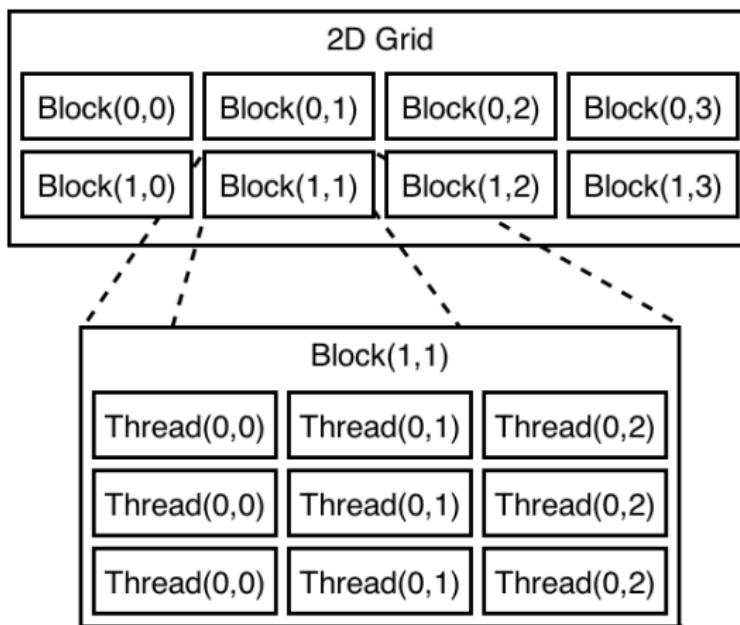
# A Possible Solution

- Think
- Wait for a fork to become available and pick up that fork
- Wait for and pick up the other fork
- Eat
- Put down the forks one by one
- Go back to thinking

# Deadlock

- Stuck when everyone holds one fork
- This is a deadlock
- Can be solved with a mediator

# Blocks



# Blocks

